

# Dynamic Statistical Profiling of Communication Activity in Distributed Applications

*Jeffrey S. Vetter*

This article was submitted to  
Association of Computing Machinery SIGMETRICS 2000  
International Conference on Measurement and Modeling of Computer  
Systems; Marina Del Rey, California; June 15-19, 2002

U.S. Department of Energy

Lawrence  
Livermore  
National  
Laboratory

**October 12, 2001**

## DISCLAIMER

This document was prepared as an account of work sponsored by an agency of the United States Government. Neither the United States Government nor the University of California nor any of their employees, makes any warranty, express or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise, does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States Government or the University of California. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States Government or the University of California, and shall not be used for advertising or product endorsement purposes.

This is a preprint of a paper intended for publication in a journal or proceedings. Since changes may be made before publication, this preprint is made available with the understanding that it will not be cited or reproduced without the permission of the author.

This report has been reproduced  
directly from the best available copy.

Available to DOE and DOE contractors from the  
Office of Scientific and Technical Information  
P.O. Box 62, Oak Ridge, TN 37831  
Prices available from (423) 576-8401  
<http://apollo.osti.gov/bridge/>

Available to the public from the  
National Technical Information Service  
U.S. Department of Commerce  
5285 Port Royal Rd.,  
Springfield, VA 22161  
<http://www.ntis.gov/>

OR

Lawrence Livermore National Laboratory  
Technical Information Department's Digital Library  
<http://www.llnl.gov/tid/Library.html>

# Dynamic Statistical Profiling of Communication Activity in Distributed Applications

Jeffrey S. Vetter

Center for Applied Scientific Computing  
Lawrence Livermore National Laboratory  
vetter3@llnl.gov

---

A complete trace of communication activity for a terascale application is overwhelming in terms of overhead and storage. We propose a novel alternative that enables profiling of the application's communication activity using statistical message sampling during runtime. We have implemented an operational prototype and our evidence shows that this new technique can provide an accurate, low-overhead, tractable alternative for performance analysis of communication activity. Moreover, this alternative enables an assortment of runtime analysis techniques not previously available with post-mortem, trace-based systems. Our assessment of relative performance and coverage of different sampling and analysis methods shows that purely random selection is preferred over counter- and timer-based sampling. Experiments on several applications running up to 128 processors demonstrate the viability of this approach. In particular, on one application, statistical profiling results contradict conclusions based on evidence from tracing. The design of our prototype reveals that parsimonious modifications to the MPI runtime system could facilitate such techniques on production computing systems, and it suggests that this sampling technique could execute continuously for long-running applications.

---

## 1 Introduction

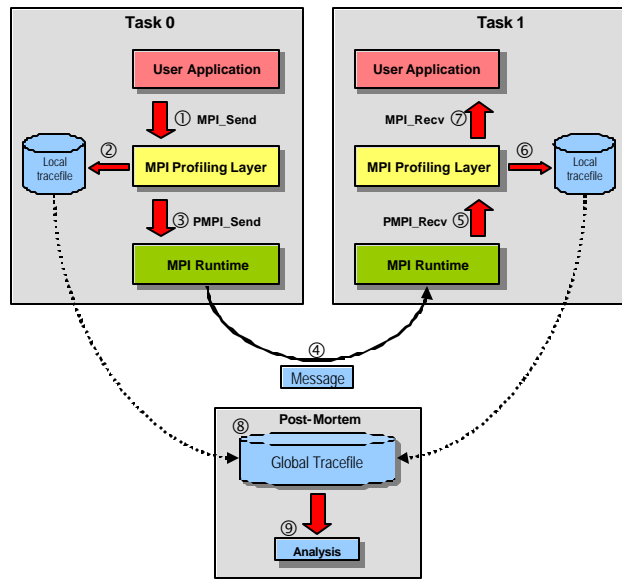
To fully realize the potential of terascale computing, users must be able to understand the performance of their applications. Unfortunately, the scale of new systems, which will have thousands, if not millions, of processors [1], is quickly outstripping the capabilities of traditional performance analysis techniques. While traditional trace-based techniques for analyzing communication performance of distributed applications have demonstrated advantages [6, 10, 11, 14, 15, 17, 18, 21], their operation on terascale platforms presents several challenges. In particular, these techniques require post-mortem analysis of potentially massive tracefiles, which, in turn, can lead to high instrumentation overhead and perturb performance observations.

Put simply, this paper proposes a novel alternative that addresses these challenges by enabling statistical profiling for individual messages of an application's communication activity during execution. Similar to other statistical profiling techniques [2, 3, 7], our technique strikes a balance between the comprehensive detail of tracing and the insight necessary for optimization. Evidence from an operational prototype, built on the Message Passing Interface (MPI), shows that this new technique can provide an accurate, low-overhead, tractable alternative for performance analysis. Messages can be sampled and analyzed with a variety of techniques that are easily interchanged at the MPI profiling layer. Also, this alternative enables an assortment of runtime analysis techniques not previously available with post-

mortem techniques, allowing the prototype to jettison raw performance data as soon as the runtime analysis is complete.

## 1.1 Motivating Example

To motivate the demands of performance analysis with large-scale applications, we consider a case study of SMG2000, a MPI application that has demonstrated scalability to four thousand processors. (Section 4 provides complete details of the experimental evaluation.) The goal of this example is to outline the process of traditional performance analysis, highlight its limitations, and argue for statistical profiling of communication activity via message sampling.



**Figure 1: Traditional Performance Analysis of Communication Activity.**

Trace-based performance analysis of distributed applications is very useful because it provides users with detailed chronology of their application's execution [6, 11, 14, 15, 17, 18, 21]. The typical operation of a trace-based tool for analyzing communication operations on a distributed application is a multi-step process as illustrated in Figure 1. Assume that the user has instrumented their application with software instrumentation that captures pertinent performance data (achieved by the *MPI Profiling*

*Layer* in Figure 1). As the instrumented application executes on a distributed platform, the instrumentation captures significant information (steps ② and ⑥) as an event record that includes a timestamp, subroutine parameters, and message size. Although each task stores its performance data locally on each CPU to limit perturbation to the network, the application does incur some level of mandatory overhead from this software instrumentation. At application termination, the tool merges the distributed files into one file, reconciling point-to-point communication operations and matching collective operations (step ⑦). Clearly, certain types of performance analysis cannot occur until this phase because metrics, such as message latency are not available until reconciliation of the event records for

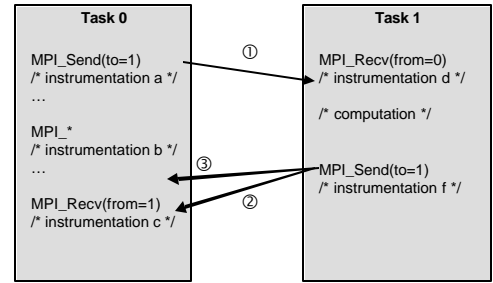
point-to-point messages. Lastly, with the merged file in hand, a user can proceed with the investigation with a variety of techniques (step ⑨) including statistical analysis, pattern recognition, automated

classifiers, and visualization to glean important insights into their application performance. Practically all of the popular tools in this area rely on visualization to assist users in analysis of performance data [11, 20].

## 1.2 Observations

To make this process more concrete, we apply a widely used MPI tracing tool to SMG2000 on 48 tasks. This application sets up and solves a linear system, a task common to scientific computing. Our tracing tool interposes instrumentation between the application and the MPI runtime using the MPI profiling layer. This instrumentation intercepts significant MPI subroutines, captures a timestamp and relevant subroutine parameters, records the event to a memory buffer, and eventually, writes them to a local disk file. Most of this instrumentation has been optimized to use efficient buffering techniques, low overhead timers, and minimal data collection. This tool automatically merges the data into one file at application termination, during which it sorts the events by global time, reconciles point-to-point message operations, and calculates communication statistics.

SMG2000 has an astonishing number of communication operations: it sends approximately 16,000 messages per task per solve. This volume of message traffic can create very large local and global tracefiles. For example, the final, merged tracefile for SMG2000 on 48 tasks is 225MB. More problematic is the execution time of SMG2000's solution phase increases from 26 seconds to 66 seconds. Consequently, it increases overhead due to the cumulative effect of the software instrumentation and the fact that the buffers must be flushed to disk frequently. Even when a user exercises considerable care in focusing the instrumentation on particular subroutines or on limited phases of execution, this situation can perturb the underlying application so much that it does not resemble an actual execution of the optimized application, especially when tracing a sequence of messaging operations. For instance, Figure 1 illustrates a sequence of communication operations that instrumentation perturbation can influence. Once cause of performance problems in MPI applications is handling unexpected messages. Typically, users try to



**Figure 2: Perturbation example.**

optimize their applications by posting receives before their matching sends as for messages ① and ②. Tracing tools would certainly provide insight into this phenomenon; however, the intervening instrumentation may very well change a properly posted message into an unexpected message. In this example, the cumulative instrumentation of  $a$  and  $b$  could very well delay the `MPI_Recv` and make message ③ appear as an unexpected message, when in reality, without instrumentation, it is not. Statistical profiling of messages helps to alleviate this issue in two ways. First, it reduces the instrumentation overhead because fewer messages are measured. Second, statistical profiling randomly distributes the instrumentation overhead across the entire message population, helping to avoid situations like the one proposed in Figure 2.

## 2 Enabling Runtime Analysis of Communication Activity

As identified in Section 1.2, one major limitation of current techniques is that most analysis must be deferred until the distributed tracefiles can be merged and reconciled. In contrast to these trace-based techniques, we propose a new alternative that enables runtime analysis of communication activity by appending a tiny amount of performance information to all messages exchanged by the application. Our prototype shows that these parsimonious modifications enable runtime analysis and statistical sampling of messages, both of which are important elements in eliminating the burdens incumbent on trace-based techniques: data management and overhead.

PHOTON, our operational prototype for statistical sampling of application communication activity, focuses on the Message Passing Interface (MPI) [8, 19]. Historically, users have written scientific applications for large distributed memory computers using explicit communication as the programming model. This trend crystallized with the creation of the Message Passing Interface (MPI) specification [8, 19], which simplified numerous issues for both application developers and system designers. As a result, application developers stabilized on the MPI programming model and this has facilitated the ongoing development of a considerable number of applications based on MPI. MPI provides a wide variety of communication operations including point-to-point operations, both blocking and non-blocking, and collective operations such as broadcast and global reductions. We concentrate on basic point-to-point

operations: blocking send, blocking receive, non-blocking send, and non-blocking receive, because they are widely used and the most important yet difficult components to sample.

## 2.1 Components

PHOTON's design has two basic components: a MPI profiling library and a modified MPI runtime library. A user application must use both components to benefit from PHOTON features. This design minimizes modifications to the underlying MPI runtime, while retaining considerable flexibility for selecting sampling and analysis techniques in the MPI profiling layer. Most MPI performance analysis techniques use only the MPI profiling layer to gather performance information.

As Figure 3 illustrates, the first component of PHOTON is the modified MPI runtime library. Our implementation is a fully functional version of MPICH 1.2.2, configured to use IBM's Message Passing Library (MPL). Our modifications are minimal.

We change the definition of point-to-

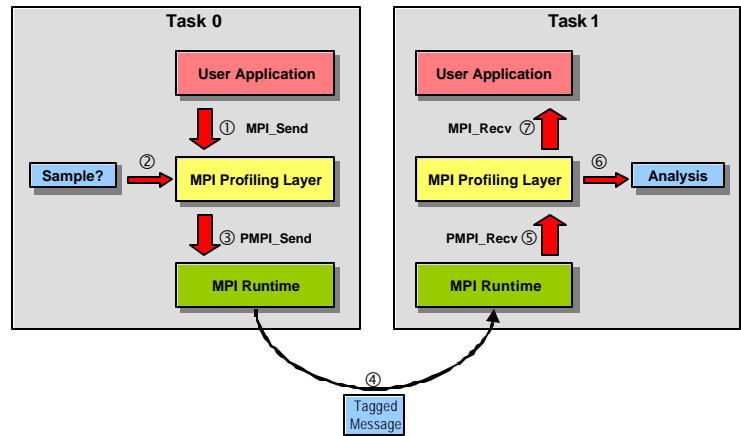


Figure 3: PHOTON Design Overview.

point message headers to include two new fields: a timestamp and a source code location identifier. We also modified the send and receive operations for all message protocols. For send operations, the profiling layer sets two task variables. When set, the modified send operations copy these two variables into the header of the outgoing message. For receive operations, the internal library copies these two variables into a modified `MPI_Status` variable. This new definition of `MPI_Status` includes our two new fields, since each incoming message must set a status word<sup>1</sup>. If the incoming message is tagged, then the receive operation copies the timestamp and source code location into these two status fields. Otherwise, it ignores them. The profiling layer can then easily check these two new fields in the `MPI_Status` structure to determine if, indeed, the send operation tagged the current message. If it did, then the receive operation, can extract these two additional fields from the `MPI_Status` structure, and use them for analysis.

Strictly speaking, our limited changes to the MPI runtime system are 1) increased message header size by 12 bytes, 2) two writes to these fields in the send operation, 3) two reads of these fields in the receive operation, 4) one control branch each in the send and the receive operation, and 4) an increased size of the `MPI_Status` structure. As our experimental results show in the next section, this overhead is imperceptible in the performance of our MPI implementation. Notice that these changes do not include procedure calls, data analysis, or sampling logic. The *optional* performance analysis tool in the MPI profiling layer provides all of these components, if desired.

The second component of our strategy exploits this additional information by using the MPI profiling layer to allow analysis on unmodified application codes. This relinking via the profiling layer interposes PHOTON between the application and the MPI runtime system, where PHOTON can intercept all MPI calls and record sufficient information about these operations to make reasonable judgments about the application performance.

## 2.2 Operation

Given these two components, PHOTON works as illustrated in Figure 3. A message is sent from Task 0 to Task 1 using MPI's blocking communication routines. In ①, the application prepares the message and calls the `MPI_Send` routine. The `MPI_Send` routine is intercepted by the MPI profiling layer. At this point ②, PHOTON uses some technique to decide whether to sample this particular message. Assume that it decides to sample the current message. PHOTON then records the start timestamp of the send operation and an identifier describing this operation's location in the source code. Usually, this location is the return address of this MPI call; however, it can be a more elaborate hash function that encodes a stack traceback, message parameters, etc. PHOTON passes these two additional pieces of information to the MPI runtime and calls the name-shifted profiling layer routine, `PMPI_Send`. As the MPI runtime begins at step ③, it loads this additional information into the message header of the message envelope that it prepares. It then dispatches this message to the underlying message libraries. This operation is similar for all MPI protocols. As the tagged message flows across the network at step ④, it carries these two additional pieces of information with it in its message header.

---

<sup>1</sup> The MPI specification allows users to set `MPI_STATUS_IGNORE` to bypass the setup of the `MPI_Status` structure.



Meanwhile, Task 1 has issued a blocking `MPI_Recv` for this message, not knowing if the incoming message is a tagged message or not. As the user application calls `MPI_Recv` at step ⑦, it must post the blocking `PMPI_Recv` at step ⑤ without knowledge of the sampling decision. Only when the message is received can Task 1 actually make a decision about how to handle this message. It is important to note that this problem is impossible to solve for all MPI protocols within the MPI profiling layer alone, and it necessitated our modifications to the MPI runtime.

Now, Task 1 receives the message at step ④ and it recognizes that this message has been tagged, so the MPI runtime copies this information directly into the two additional fields in our modified `MPI_Status` structure. When the MPI runtime completes the receive of this tagged message, it returns from the `PMPI_Recv` into the PHOTON profiling layer. At this point, PHOTON can make any number of decisions about how to analyze the tagged message. At ⑥, PHOTON can record statistics and discard the data, write it to a trace file, or simply ignore it. When PHOTON has completed its analysis at step ⑥, it returns to the user application via the `MPI_Recv` call. The user application can then process the message as it normally would and without regard to the fact that the message was sampled by the underlying performance analysis system.

## 2.3 Key Design Implications

This design alternative has a number of important implications. First, message processing is undisturbed. Our technique does not require additional messages, additional copying of message buffers, or excessive quantities of extra buffer space. Therefore, the MPI behaviors of this alternative should closely resemble the behaviors of the original MPI application. Although we considered several options to modifying the underlying MPI implementation, none of these strategies provided necessary functionality for all MPI operations and respected all of the requirements demanded by the MPI specification, such as its message-ordering requirement.

More specifically, three alternatives come to mind. In the first alternative, one could simply send an extra message following each sampled message. This alternative has two issues: it introduces additional messages into the system, and the receiving task has no *a priori* knowledge of when to wait for a sampled

---

PHOTON creates a temporary `MPI_Status` structure and passes this structure to the underlying system.

message. This strategy could also introduce race conditions into the application. In the second alternative, the system could exchange performance data during collective operations; however, the overhead could be noticeable and reconciling sends with receives would still demand comprehensive knowledge of all operations. The final alternative is that the technique could use MPI derived types to piggyback additional data onto messages. On the face of it, this alternative is appealing. The MPI specification allows nested data types, allowing a send operation to simply repackage a message and forward it to the receiver where it is unpackaged with the additional performance data. Unfortunately, this alternative has several problems. First, MPI derived types can perform poorly because of additional memory copying and buffering of data, which might drastically alter the performance characteristics of the application [9]. Second, the receiver cannot determine which messages are tagged messages. MPI message envelopes specify message source, tag, and communicator, but not data type. This limited information prohibits the receiver from determining whether an incoming message is tagged. Therefore, the receive operation has no idea of how to prepare the receive buffer. This strategy is also plagued by the possibility of several types of race conditions.

The second important implication is that decisions regarding sampling and analysis techniques remain at the MPI profiling layer. Thus, these decisions can be easily interchanged without altering the underlying MPI runtime. Better still, by relegating all of these decisions to the profiling layer, the performance of the modified MPI runtime can remain practically unchanged from the original.

The third and final implication is that this design is applicable to all types of point-to-point communication regardless of MPI subroutine or message protocol. Take, for example, non-blocking communication. PHOTON simply loads the message header during the initiation of the non-blocking `MPI_Isend`. The underlying message library transfers the message normally. As the message is received with a `MPI_Irecv/MPI_Wait` pair, the extra performance information from the message header is transferred directly into the `MPI_Status` structure, which is provided to `MPI_Wait` as a parameter. This line of reasoning holds true for other completion operations including `MPI_Test` and `MPI_Waitsome` too because they return a `MPI_Status` structures for completions.

### 3 Statistical Message Sampling

Although sampling is popular in many other areas of performance analysis, such as procedure profiling [3, 7] and instruction analysis using hardware counters [2], it has not been applied to communication activity due to the limitations listed in Section 1.2. With the novel alternative proposed in Section 2, we can now reliably and accurately access performance information at runtime, so that communication performance analysis could benefit equally from statistical sampling. To our knowledge, this technique is novel and it represents a significant shift in current technology for performance analysis of terascale, distributed applications. Sampling has also been applied at low levels of communication activity [5], but this research focused on understanding wide-area networks, rather than on optimizing the applications that use those networks.

#### 3.1 Sampling Strategies

As Figure 3 illustrates, in step 2, PHOTON can use most any technique to determine how to sample messages from the entire message population of the application. Naturally, three different techniques can drive our sampling approach: purely random sampling, counter-based sampling, and timer-based sampling.

Random sampling: Our first sampling method is purely random sampling. On every send operation, PHOTON draws a number from a uniform distribution in  $(0,1]$  and then checks that number against a user defined threshold ( $T$ ) to determine if the current message should be sampled. This strategy is simple and it allows a user to easily control the number of samples by changing the threshold, and the instrumentation impact on the application.

Counter-based sampling: In PHOTON, a single counter in each task is incremented for every send operation. When the counter exceeds a threshold, one message is sampled, the counter is reset, and a new target threshold is calculated. The user can select the period ( $P$ ) and variance ( $V$ ) of the counter. Counter-based sampling benefits from its simplicity and low overhead; however, estimating the appropriate settings for  $P$  and  $V$  can be difficult because it depends entirely on the frequency of application communication.

Timer-based sampling: Similarly, for timer-based sampling, a message is sampled after a period of time has expired since the last send operation. The user specifies a period of time (P) and a variance (V). Note that our technique does not use expensive interrupts to execute this sampling technique. Rather, as the application calls MPI send routines, PHOTON polls the local time and then decides if the specified threshold has been met. Although the cost of sampling the timer can be expensive relative to using counters, timer parameters are much easier to estimate and PHOTON must capture this timestamp for the outgoing message.

These methods sample from the entire message population; however, they can also be adapted to sample subsets of the population. For instance, we could sample only large messages, only messages sent from a certain callsite, or only messages during a certain phase of the application execution.

### **3.2 Analysis Methods**

Now that a considerable amount of performance information is available at runtime, PHOTON can elect to perform analysis at runtime and jettison the raw performance data. In the context of terascale computing, this capability is vital because it eliminates the need for capturing massive tracefiles and harvesting important performance data from those files. Certain performance problems may still require tracing; however, this runtime analysis can quickly direct further efforts on a subset of operations.

For performance analysis, our experience indicates that it is important that users are able to map performance data back to source code; we record and categorize all messages by source task, destination task, callsite location in source and destination, and message size. We capture message latency as our primary performance metric, where we define latency as the time from the start of the send operation until the end of the matching receive operation. Although this definition is somewhat different than architectural definitions, this interpretation is easy for a user to reason about, and it maps directly to the user's source code. Most trace-based tools use similar definitions.

We introduce two lightweight techniques to analyze this data at runtime and we include a expensive technique of writing the event to a local file, as a reference point.

Statistical summary (STAT): One traditional option for profiling is a statistical summary of the messages tabulated by source code locations of the send and receive operations. For our statistical

summary, PHOTON captures a maximum, a minimum, a count, and a cumulative total of the message latency. With a reasonably long application execution, this statistical information would supply a convincingly accurate picture of communication activity. It also has a very small analysis overhead that includes locating an entry in a data structure and updating several fields in that entry. This summary when combined with the message size and topology information can identify performance anomalies of individual message operations.

Frequency distributions (FREQ): Another valuable technique for analyzing large masses of raw data is a frequency distribution. Using a frequency distribution, an analyst can project any number of raw data samples into classes, which are easy to represent, understand, and store. These distributions provide more information than the statistical summary presented above; however, the overhead for updating the entry and the memory requirements are slightly increased. As before, we categorize the data by sender, receiver, callsite location for both, and message size. Our implementation creates an array of buckets for message latency delimited by a common log scale. Each bucket counter is incremented when a message's latency falls within that bucket's bounds.

Write to File (WRITE): Although our main focus for PHOTON is runtime analysis of performance data, we can simply use the analysis step write the values to a memory buffer or a file. In contrast to traditional tracing, however, this file would not represent a complete chronology of the application's communication behavior, but only a small, sampled portion of the overall communication.

## 4 Evaluation

Our evaluation focuses on our hypothesis that PHOTON can offer a feasible alternative to trace-based analysis of message passing applications on terascale platforms. In this regard, we evaluate PHOTON along several dimensions: overhead and perturbation, sampling and analysis methods, and improvements in data management. Then, we apply PHOTON to several applications in a realistic situation.

### 4.1 Platform

We ran our tests on an IBM SP system. This machine is composed of sixteen 222 MHz IBM Power3 8-way SMP nodes, totaling 128 CPUs. Each processor has three integer units, two floating-point units, and two load/store units. Its 64 KB L1 cache is 128-way associative with 32 byte cache lines and L1 uses

a round-robin replacement scheme. The L2 cache is 8 MB in size, which is four-way set associative with its own private cache bus. At the time of our tests, the batch partition had 15 nodes and the operating system was AIX 4.3.3. Each SMP node contains 4GB main memory for a total of 64 GB system memory. A Colony switch--a proprietary IBM interconnect--connects the nodes. We compiled the various tests with the IBM XL compilers. Our test jobs ran on dedicated nodes, although other jobs were concurrently using the network. We built our prototype of PHOTON with the publicly available version of MPI from Argonne National Laboratory: MPICH 1.2.2. We configured MPICH to use IBM's Message Passing Layer (MPL) as the communication substrate.

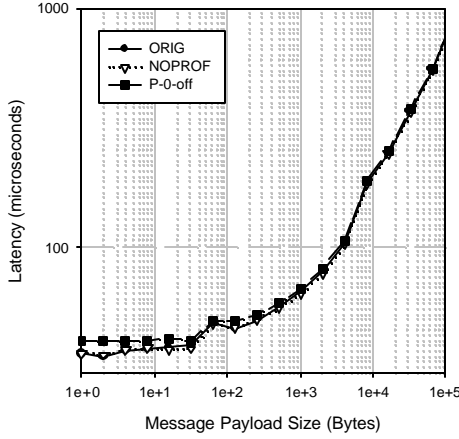
Configuration	Moniker	PHOTON MPI Runtime	PHOTON profiling layer	Sampling method	Analysis Method	Description
Original MPI	ORIG	No	No	N/A	N/A	
MPI Modified to include PP information	NOPROF	Yes	No	N/A	N/A	MPI Runtime w/ PHOTON modifications. Lacks PP profiling layer that includes sampling and analysis methods.
Sampling enabled at various random thresholds	P-X	Yes	Yes	Random X = threshold	None	Modified MPI with sampling infrastructure installed where the threshold for random sampling is X. No analysis is performed on the sampled message.
Sampling enabled at various counter rates	C-X	Yes	Yes	Counter X = period	None	Modified MPI with sampling infrastructure installed where the counter-sampling period is X. No analysis is performed on the sampled message.
Sampling enabled at various timer periods	T-X	Yes	Yes	Timer X = period	None	Modified MPI with sampling infrastructure installed where the timer-sampling period is X. No analysis is performed on the sampled message.
Sampling w/ various analysis methods	A-X-B	Yes	Yes	A = sampling method X = threshold	B is print, stat, or freq	Modified MPI with sampling infrastructure installed where the sampling method is A with threshold X and the stated analysis B is performed on each sampled message.

**Table 1: Configuration Overview.**

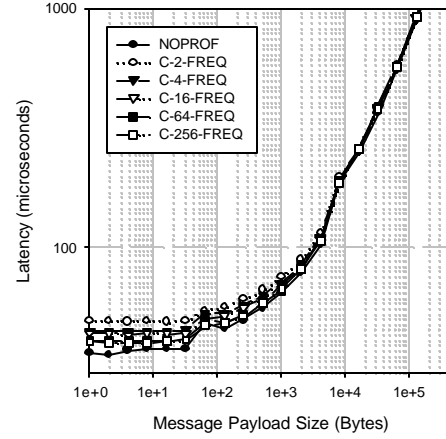
## 4.2 Overhead

Our proposals hinge on the ability of the MPI runtime system to carry a tiny amount of additional information with each message. To assess the penalty for this increase in the message header size, we compare three versions of MPI as described by Table 1: 1) a version of the original MPI (ORIG); 2) a version of the modified MPI the minimal internal changes to the internals, but without the PHOTON profiling layer installed (NOPROF); and, 3) a version with both the modified MPI and the PHOTON

profiling layer configured to sample no messages (P-0-OFF). For these experiments, we did not perform analysis on the resulting sampled message; we consider them in the next section. These experiments measure message latency for each configuration using the Pallas PMB-MPI1 PingPong benchmark between two SMP nodes across the network.



**Figure 4: PHOTON Overhead.**



**Figure 5: PHOTON Overhead for Counter Sampling.**

As Figure 4 illustrates, the modified version of MPI (NOPROF) performs identically to the original version (ORIG). The increase in latency was less than one microsecond for all message sizes. When we included the PHOTON profiling library with sampling disabled (P-0-OFF), the overhead increased by 3 microseconds at a payload size of 4 bytes, though it is only noticeable for smaller messages that have fewer than 10,000 bytes. As the message size increases, however, this fixed cost overhead disappears into the greater cost of the communication operation, as expected. Although we do expect the overhead for PHOTON to be very small for most MPI implementations, the impact of these changes will proportionally higher on optimized messaging layers, such as those that use shared memory for intra-node transfers on an SMP. Our experiments showed no measurable change in the communication bandwidth across these configurations.

This result is pivotal because it argues that these new features could be included in many MPI implementations without mandatory performance degradation. The optional profiling library does inflict a small overhead, but it is only noticeable only for small messages.

### 4.3 Sampling Methods

The various sampling methods introduced in Section 3.1 have different performance penalties and dissimilar results with respect to how they sample the message population of the application. All three options provide a convenient means to control the number of messages sampled, and we use that feature to evaluate each method's overhead and sample space. We delay discussion of sample space until Section 0 where we evaluate it on real applications. As shown in Table 1, we vary the threshold (X) for random sampling (P-X), the period (X) for counter-based sampling (C-X), and the period (X) for timer-based sampling (T-X). Each sampled message is analyzed with the frequency distribution technique.

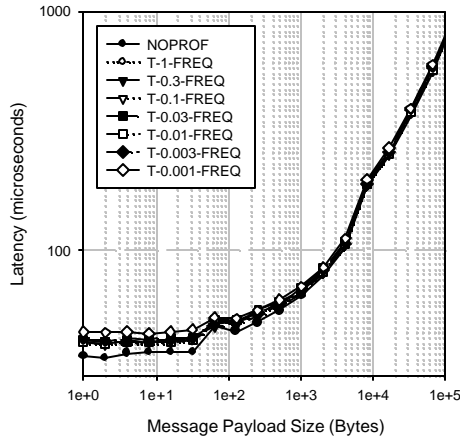


Figure 6: PHOTON Overhead for Timer Sampling.

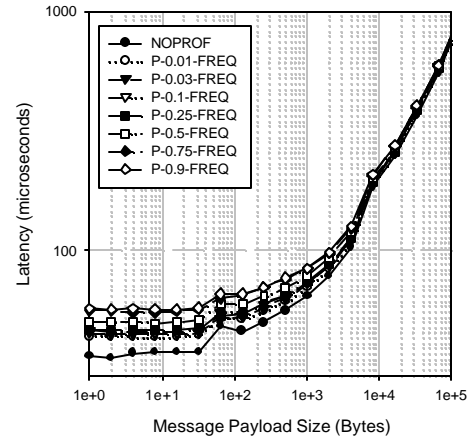


Figure 7: PHOTON Overhead for Random Sampling.

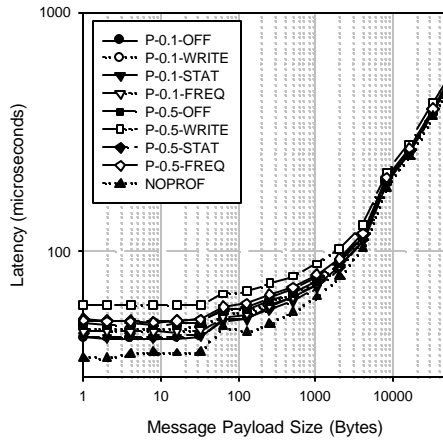
First, for counter sampling, we range the period from 2 to 256. Figure 5 illustrates the overhead's explicit variation with the counter period. For a payload of 4 bytes, the latency ranges from 48.7 microseconds for C-2 down to 40.1 microseconds for C-256. For reference, we include the lower bound latency of 37.2 microseconds at 4 bytes for the NOPROF configuration. For applications that send messages with payloads of more than 20k bytes, the cost is imperceptible relative to the overall messaging cost.

Next, Figure 6 shows that the overhead for timer-based sampling is less controllable, but it remains low. A 4byte message suffers an increase in latency from 37.2 microseconds for NOPROF to 45.5 microseconds for T-0.001-FREQ, when messages are sampled at a period of 1 millisecond. Increasing the sampling period to 300 milliseconds lowers the latency to 41.8 microseconds.



Last, for random sampling, we range the threshold from 0.01 (1%) up to 0.9 (90%). In Figure 7 illustrates the overhead's explicit variation with this threshold range. At 1% for a 4-byte message, the latency increases to 43.3 microseconds from 37.2 microseconds for NOPROF. As the threshold increases up to 90%, the latency increases to 57.0 microseconds.

This evidence reveals that all three sampling strategies have a demonstrated variance of overhead with respect to message latency.



**Figure 8: PHOTON Overhead for Analysis Methods.**

#### 4.4 Analysis Methods

Runtime analysis methods are valuable because after the analysis, the raw data can be discarded, eliminating many of the problems with traditional techniques. However, the analysis methods must balance multiple conflicting goals of minimizing computation and storage against providing enough detailed information for users to make proper attribution of performance data. For

this reason, we evaluate two lightweight statistical analysis techniques and a traditional technique of writing the data to a local file. Viewed in this light, it is important to remember that the cost of the analysis technique can be balanced against the sampling rates evaluated earlier. If a particular analysis method has relatively high cost, then the sampling rate can be lowered to compensate, minimizing the overall impact on the application execution.

As Figure 8 illustrates, the overhead of analysis methods varies considerably; however, our ability to change the sampling rates at runtime provides considerable flexibility in balancing this analysis overhead with the sampling rate. One important feature of PHOTON is that the performance data can be interpreted at runtime, and discarded when the analysis is complete. In Figure 8, we see that the most expensive technique is writing a record to a file (P-0.5-WRITE) at 59.9 microseconds, and as we would expect, that the lowest overhead occurs when analysis is disabled (P-0.1-OFF). This overhead drops as we decrease the random sampling threshold to 10% (P-0.1) from 50% (P-0.5) for each analysis technique. Although

both the statistical summary method and the frequency distribution method yield similar levels of overhead, they perform much better than the write technique.

## 4.5 Applications

To verify that our new technique works with real applications, we tested three applications: sPPM [16], Sweep3d [13], and SMG2000 [4]. All of these benchmarks have scaled to thousands of tasks. In these examples, we scaled problem sizes with the number of tasks, keeping the amount of work per task constant.

sPPM: sPPM [16] solves a 3-D gas dynamics problem on a uniform Cartesian mesh, using a simplified version of the Piecewise Parabolic Method. The algorithm makes use of a split scheme of X, Y, and Z Lagrangian and remap steps, which are computed as three separate sweeps through the mesh per timestep. Message passing provides updates to ghost cells from neighboring domains three times per timestep.

SMG2000: SMG2000 [4] is a parallel semicoarsening multigrid solver for the linear systems arising from finite difference, finite volume, or finite element discretizations of the diffusion equation  $\nabla \cdot (D \nabla u) + \mathbf{s}u = f$  on logically rectangular grids. The code solves both 2-D and 3-D problems with discretization stencils of up to 9-point in 2-D and up to 27-point in 3-D.

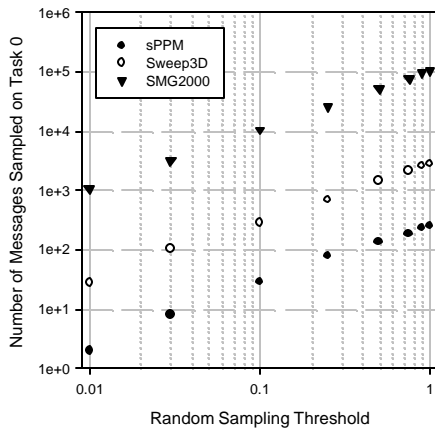


Figure 9: Number of Messages Sampled on Task 0.

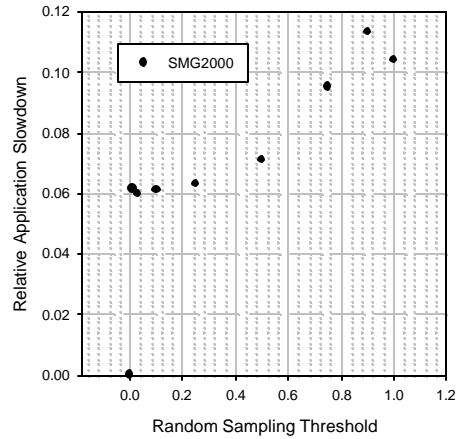
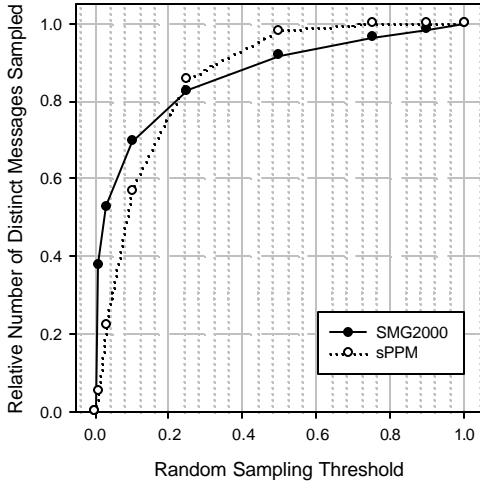


Figure 10: Impact of sampling rate on SMG2000 runtime.

Sweep3D: Sweep3D [12, 13] is a solver for the 3-D, time-independent, particle transport equation on an orthogonal mesh and it uses a multidimensional wavefront algorithm for "discrete ordinates" deterministic particle transport simulation. Sweep3D benefits from multiple wavefronts in multiple

dimensions, which are partitioned and pipelined on a distributed memory system. The three dimensional space is decomposed onto a two-dimensional orthogonal mesh, where each processor is assigned one columnar domain. Sweep3D exchanges messages between processors as wavefronts propagate diagonally across this 3-D space in eight directions.

Figure 9 shows the number of messages actually sampled by PHOTON on our three applications, where a threshold of 1.0 shows the total number of messages sent from task 0. Here we can also see the differences in the number of messages sent by every application. SMG2000 sends 105,000 messages whereas sPPM sends only 260. Encouragingly, Figure 9 illustrates that PHOTON can easily change the number of messages sampled during an application experiment.



**Figure 11: Sample Space of Message Population**

Sweep3D and sPPM send relatively few messages; our experiments confirmed that all sampling rates less than 50% did not have a noticeable effect on the application runtime. When the sampling threshold increased above 50%, the execution time of both applications increased consistently but never more than 1%. SMG2000, on the other hand, sends considerable more messages. Figure 10 shows the impact of the PHOTON sampling threshold on the SMG2000 runtime, normalized to the NOPROF configuration. We see that initially, even at a 1% sampling threshold, SMG2000 experiences a slowdown of 6.1%. As the sampling threshold increases to 90%, the slowdown climbs to 11%.

Although SMG2000 does incur perturbation on the application runtime, its impact is much less than the traditional tracing approach that had a 154% slowdown in execution time. Much more importantly,

because PHOTON gathers its information with sampling, it helps to avoid pathological measurement problems like those mentioned in Section 1.2.

Figure 12 shows a subset of the performance data from our experiments on SMG2000. The *message identifier* is a distinct identifier hashed from the tuple of (source rank, destination rank, source code location, destination code location, message size). *Latency distribution* is a frequency distribution of the message latency where bucket bounds are delimited by a common log scale.

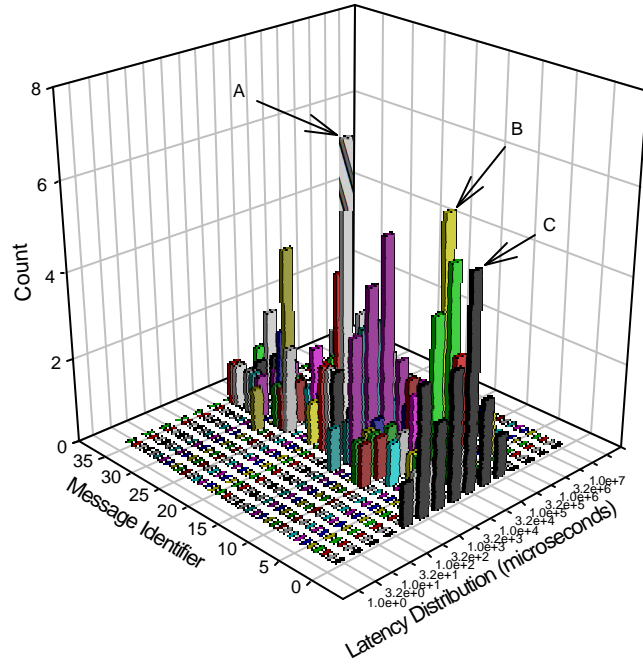


Figure 12: Example Frequency Distribution.

#### 4.6 Observations

Our experimental evaluation substantiated many of our claims and it also revealed several issues. Among the most important observations exposed by our experiments were that we could significantly control the amount of overhead on an application, statistical message sampling reduces the possibility of pathological perturbations from instrumentation, our modifications to MPI to enable runtime analysis were very small and inflicted no major performance penalty, and statistical profiling can be used on large, long-running applications for performance analysis of communication activity.

Reduced overhead. Our evaluation revealed that we could dramatically control the overhead of PHOTON and maintain important information like message latency using statistical profiling. Our experiments on real applications and on benchmarks demonstrated that this technique radically diminishes the impact of instrumentation on application execution time and message latency. Consequently, these variable sampling rates can be balanced against the cost of the analysis technique to compensate.

Parsimonious modifications to MPI have negligible performance implications. Changes to the MPI runtime system were frugal and our evidence indicates that these changes have an almost imperceptible

performance penalty for the implementation of PHOTON as the evidence of Section 4.2 demonstrates. Control of the sampling and analysis techniques remains in the profiling layer, so that users can simply interchange them with their applications to do performance analysis. Indeed, we believe this result argues for PHOTON features in production MPI implementations on terascale systems because it will be necessary for any performance analysis of communication activity. What is more, with its low overhead and amortized perturbation, this technique is also appealing to consider using continuously for large-scale, long-running applications.

Reduced data volume. Another goal of this work was to ease the data management burden imposed by traditional techniques; our experiments on several applications demonstrate that significant, useful performance information can be processed dynamically. Because we can calculate message latency immediately during runtime, the analysis can proceed immediately, and then discard the raw performance data. Our lightweight analysis examples, statistical summary and frequency distribution, showed that PHOTON can collect important information from long-running, terascale applications with little overhead.

Timer-based sampling. We find that our timer-based sampling method can lead to pathological situations where only certain messages are sampled. For example, it is common in scientific simulations to exchange data with neighbors during one phase of a timestep and compute during the other phase. We find that many of our experiments only sample the first communication operation immediately after the compute phase, and then miss the remaining send operations that immediately follow it. Our other two sampling strategies do not suffer from this deficiency because the variance allows the sampling to break out of patterns. More importantly, our random sampling strategy allows a user to simply specify a percentage of messages to sample, and PHOTON randomly selects messages without any of these concerns.

## 5 Conclusions

As parallel computing systems continue to scale to massive numbers of processors, performance analysis techniques must allow users to understand their application's behaviors. We have proposed and evaluated a novel alternative to trace-based, post-mortem performance analysis of communication activity in distributed applications. Our alternative for statistical profiling of communication activity included

minor modifications to the MPI runtime to enable runtime analysis of performance data and an implementation of message sampling. Our operational prototype, PHOTON, demonstrated several significant advantages over the traditional trace-based approaches including dramatically lower overhead and perturbation, even for applications that communicate frequently.

## Acknowledgments

This work was performed under the auspices of the U.S. Department of Energy by the University of California, Lawrence Livermore National Laboratory under contract No. W-7405-Eng-48.

## References

- [1] G.S. Almasi, C. Cascaval *et al.*, "Demonstrating the scalability of a molecular dynamics application on a Petaflop computer," Proc. Int'l Conf. Supercomputing, 2001, pp. 393-406.
- [2] J.M. Anderson, L.M. Berc *et al.*, "Continuous profiling: where have all the cycles gone?," *ACM Trans. Computer Systems*, 15(4):357-90, 1997.
- [3] T.E. Anderson and E.D. Lazowska, "Quartz: A Tool for Tuning Parallel Program Performance," Proc. 1990 SIGMETRICS Conf. Measurement and Modeling Computer Systems, 1990, pp. 115-25.
- [4] P.N. Brown, R.D. Falgout, and J.E. Jones, "Semicoarsening multigrid on distributed memory machines," *SIAM Journal on Scientific Computing*, 21(5):1823-34, 2000.
- [5] K.C. Claffy, G.C. Polyzos, and H.-W. Braun, "Application of sampling methodologies to network traffic characterization," Proc. SIGCOMM: Communications architectures, protocols and applications, 1993, pp. 194-203.
- [6] G.A. Geist, M.T. Heath *et al.*, "A Users' Guide to PICL - A Portable Instrumented Communication Library," Oak Ridge National Laboratory, P.O.Box 2009, Bldg. 9207-A, Oak Ridge, TN 37831-8083 1991.
- [7] S.L. Graham, P.B. Kessler, and M.K. McKusick, "Gprof: A Call Graph Execution Profiler," *SIGPLAN Notices (SIGPLAN '82 Symp. Compiler Construction)*, 17(6):120-6, 1982.
- [8] W. Gropp, E. Lusk, and A. Skjellum, *Using MPI: portable parallel programming with the message-passing interface*, 2nd ed. Cambridge, MA: MIT Press, 1999.
- [9] W.D. Gropp, E. Lusk, and D. Swider, "Improving the Performance of MPI Derived Datatypes," Proc. MPI Developers and Users Conference (MPIDC), 1999.
- [10] W. Gu, G. Eisenhauer *et al.*, "Falcon: On-line Monitoring and Steering of Parallel Programs," *Concurrency: Practice and Experience*, 10(9):699-736, 1998.
- [11] M.T. Heath, A.D. Malony, and D.T. Rover, "Parallel performance visualization: from practice to theory," *IEEE Parallel & Distributed Technology: Systems & Applications*, 3(4):44-60, 1995.
- [12] A. Hoisie, O. Lubeck *et al.*, "A General Predictive Performance Model for Wavefront Algorithms on Clusters of SMPs," Proc. ICPP 2000, 2000.
- [13] K.R. Koch, R.S. Baker, and R.E. Alcouffe, "Solution of the First-Order Form of the 3-D Discrete Ordinates Equation on a Massively Parallel Processor," *Trans. Amer. Nuc. Soc.*, 65(198), 1992.
- [14] J. Labarta, S. Girona *et al.*, "DiP: A Parallel Program Development Environment," CEPBA, Barcelona, Spain 1996.
- [15] A.D. Malony and D.A. Reed, "Visualizing Parallel Computer System Performance," in *Parallel Computer Systems: Performance Instrumentation and Visualization*, M.S. Bucher, Ed. New York: ACM, 1990.
- [16] A.A. Mirin, R.H. Cohen *et al.*, "Very High Resolution Simulation of Compressible Turbulence on the IBM-SP System," Proc. SC99: High Performance Networking and Computing Conf. (electronic publication), 1999.

- [17] D.A. Reed, P.C. Roth *et al.*, “Scalable performance analysis: the Pablo performance analysis environment,” Proc. Scalable Parallel Libraries Conf., 1994, pp. 104-13.
- [18] S. Shende, A.D. Malony *et al.*, “Portable profiling and tracing for parallel, scientific applications using C++,” Proc. SIGMETRICS Symp. Parallel and Distributed Tools (SPDT), 1998, pp. 134-45.
- [19] M. Snir, S. Otto *et al.*, Eds., *MPI--the complete reference*, 2nd ed. Cambridge, MA: MIT Press, 1998.
- [20] J. Stasko, J. Domingue *et al.*, Eds., *Software Visualization: Programming as a Multimedia Experience*,. Cambridge, MA: MIT Press, 1998.
- [21] C.E. Wu, A. Bolmarcich *et al.*, “From Trace Generation to Visualization: A Performance Framework for Distributed Parallel Systems,” Proc. SC2000: High Performance Networking and Computing, 2000.